

BACHELORARBEIT

**Implementation and Evaluation of a  
Support Vector Machine on an 8-bit  
Microcontroller**

ausgeführt zum Zwecke der Erlangung des akademischen Grades  
eines Bachelor of Science

unter der Leitung von

Univ.Ass. Dipl.-Ing. Dr.techn. Wilfried Elmenreich  
Institut für Technische Informatik  
Fakultät für Informatik  
Technische Universität Wien

durchgeführt von

Thomas Nowak  
Matr.-Nr. 0425201  
Sturzgasse 1C/10, A-1140 Wien

Wien, im Juli 2008

.....

# Implementation and Evaluation of a Support Vector Machine on an 8-bit Microcontroller

Support Vector Machines (SVMs) can be used on small microcontroller units (MCUs) for classifying sensor data. We investigate the theoretical foundations of SVMs, we prove in particular a very general version of Mercer's theorem, and review the software package  $\mu$ SVM, which is an implementation of an SVM for use on MCUs. Present SVM solutions are not applicable to MCUs, because they need too much memory space which can be expected when running on a personal computer, but not on an MCU. It is shown that, while  $\mu$ SVM's execution time does not scale very well to a large number of training examples, it is possible to prematurely terminate the training process and still retain good numerical accuracy.

# Contents

<b>1. Introduction</b>	<b>1</b>
1.1. Motivation . . . . .	2
1.2. Structure of the Thesis . . . . .	3
<b>2. Optimization and SVMs</b>	<b>4</b>
2.1. Basic Optimization Theory . . . . .	4
2.1.1. Convexity – Definition and Simple Properties . . . . .	5
2.1.2. Global/Local Minima . . . . .	8
2.1.3. Differentiability . . . . .	8
2.1.4. A Minimality Condition . . . . .	10
2.2. Bayesian Classification . . . . .	11
2.3. Constrained Quadratic Optimization for SVMs . . . . .	13
2.4. Are we there yet? – Optimality Conditions . . . . .	14
2.5. Kernel Functions . . . . .	16
<b>3. SVM Algorithms</b>	<b>19</b>
3.1. Naïve KKT . . . . .	19
3.2. Dealing with Errors . . . . .	20
3.3. The Dual Formulation . . . . .	21
3.4. Solution Strategies . . . . .	22
3.4.1. Osuna et al. . . . .	22
3.4.2. SMO . . . . .	23
3.4.3. SVM <sup>light</sup> . . . . .	25
3.5. The Classifier Function . . . . .	26
<b>4. Implementation</b>	<b>27</b>
4.1. $\mu$ SVM Overview . . . . .	27
4.2. Target Hardware . . . . .	28
<b>5. Results and Discussion</b>	<b>30</b>
5.1. Performance . . . . .	30
5.2. Speed of Convergence . . . . .	31
<b>6. Summary</b>	<b>33</b>
<b>Bibliography</b>	<b>34</b>
<b>A. Notation</b>	<b>35</b>

<b>B. <math>\mu</math>SVM Documentation</b>	<b>36</b>
B.1. Training . . . . .	36
B.1.1. Kernels . . . . .	37
B.1.2. Early Termination . . . . .	37
B.2. Classification . . . . .	37
B.3. Compiler Flags . . . . .	37
B.4. Example . . . . .	37

# 1. Introduction

Support Vector Machines (SVMs) are a method in machine learning where the goal is a binary classification of input data. The procedure is the following: In a first phase, the so-called *training phase*, a set of labeled objects is presented to the machine. The labels are taken from a two-element set (e. g., 0/1, A/B, +1/-1, good/bad, ...). The machine's objective is to derive a decision procedure from this *training set* such that it can classify correctly objects presented to it in the future. Of course, it is the supervisor's responsibility to choose the training set adequately, i. e., taking as representative objects as possible to enable good classification performance on future observations.

Support Vector Machines are used for a variety of classification tasks. These include handwriting recognition, speaker identification, text categorization and face detection.

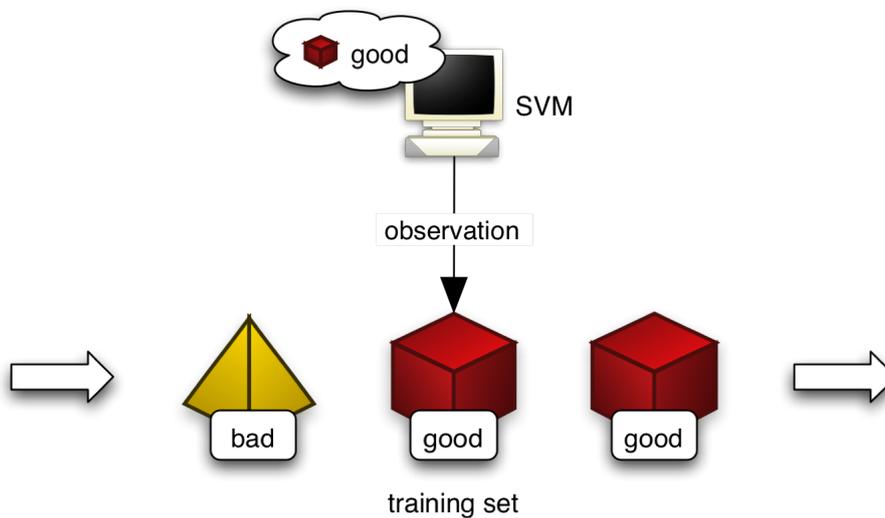


Figure 1.1.: SVM Training

## 1.1. Motivation

In a technical device, e. g., a control loop, it is necessary to have some sort of input data from the environment. These values are the output of *sensor devices*. Now, relying on a single sensor is not reasonable since the sensor may fail or behave highly differently in different environmental situations (e. g., temperature, light, movement). To compensate for this, sensors are replicated (not necessarily identically) and even completely different types of sensors may be used. Consider, for example, the velocity sensor of a roller coaster wagon. It might be known that velocity sensor A works well in the temperature range of 0°C to 20°C and that velocity sensor B is more accurate from 20°C onwards. So, it would be reasonable to install both sensors A and B and additionally put a temperature sensor on the device. In that way, one could use sensor A up to 20°C and sensor B above. The value returned by the sensor device will be that of sensor A or sensor B depending on the temperature. It is transparent to the control application which sensors are used and even which sensors exist in the device. The calculation<sup>1</sup> of the returned value is done by a microcontroller unit (MCU) inside the sensor device.

Now, it could be possible that the choice which sensor to trust is based on more than one parameter. Then the decision which sensor to use would no longer be as simple as the check whether one number is smaller than the other, but a more sophisticated method would have to be used. This problem can be tackled with the use of an SVM.

Of course, the previous example alone does not justify an implementation of the SVM *training* algorithm on the MCU, because most often, the characteristics of the sensors are known beforehand and the training can happen offline on a PC. But situations may occur where such a priori knowledge is not accessible, for example if the device is to be used in a broad variety of different environments which are not determined a priori and the device has to adjust to a new environment periodically.

Other applications for an SVM implementation on an MCU might include classification tasks like determining a “safe state” in an execution<sup>2</sup> or image classification (though a pixel-by-pixel comparison will not be feasible).

---

<sup>1</sup>It should be noted that this calculation need not be as easy as choosing the right sensor.

It might involve taking different sensor values with appropriate weighting factors as well as more complex methods. The use of such methods is referred to as *sensor fusion*.

<sup>2</sup>A *safe state* could be one where it is possible for the device to enter a sleep mode, because it will not be used for a while. Of course, this decision will be a probabilistic one and hence is not suitable for safety critical applications.

## **1.2. Structure of the Thesis**

Chapter 2 gives a textbook-like introduction to the basic concepts of the branch of optimization theory needed for SVMs. Chapter 3 derives basic algorithmic concepts in SV training and reviews popular approaches. The implementation of  $\mu$ SVM is described in chapter 4. Chapter 5 states results concerning the performance of  $\mu$ SVM. The thesis ends with a conclusion in chapter 6.

## 2. Optimization and SVMs

This chapter treats the mathematical foundations of Support Vector Machines. We start with the classification of the main problem in Support Vector training – namely solving the quadratic program that yields the separating hypersurface for our training set. We then investigate optimality conditions for the optimization problem and give first ideas for efficient algorithms. We end the chapter with a slight generalization of the problem which is of greatest importance in practice (kernel functions instead of scalar product).

The chapter is essentially self-contained although many proofs are left out in order to avoid an overly mathematical bias. Only some of the most prototypical (or short) proofs are stated.

### 2.1. Basic Optimization Theory

**Definition 1.** Let  $f : X \rightarrow \mathbb{R}$  be a function (the objective function) and  $S \subset X$  be a nonempty set (the feasibility set). The problem of finding an  $x_0 \in S$  such that

$$f(x_0) \leq f(x) \quad \text{for all } x \in S \quad (2.1)$$

is called a minimization problem.  $\square$

We note that trying to solve this problem only makes sense if  $\inf\{f(x) \mid x \in S\}$  exists, i. e.,  $f$  is bounded from below on  $S$ .

Analogously, if we exchange “ $\leq$ ” by “ $\geq$ ” in Definition 1, we get a *maximization problem*. We call both an *optimization problem*. We will only be concerned with minimization problems here as the other case is analogous – we just have to flip a few inequality signs, exchange min by max, inf by sup, etc.

The problem that we will try to solve in SVM training is a so-called *convex optimization problem*, i. e., both the objective function  $f$  and the feasibility set  $S$  are convex. We will define these terms and derive some first properties.

For the following, we set

$$\overline{\mathbb{R}} = \mathbb{R} \cup \{+\infty\} \quad (2.2)$$

with

$$x < +\infty \text{ and } x + \infty = +\infty \quad (2.3)$$

for every real  $x$ . Further,  $X$  will always denote a subset of a finite-dimensional Hilbert space  $H$  over the reals. We can think  $X \subset \mathbb{R}^n$  here.

### 2.1.1. Convexity – Definition and Simple Properties

**Definition 2.** Let  $A \subset X$ . We call  $A$  convex if for all  $x, y \in A$  and  $0 < \lambda < 1$ :

$$\lambda x + (1 - \lambda)y \in A \quad (2.4)$$

□

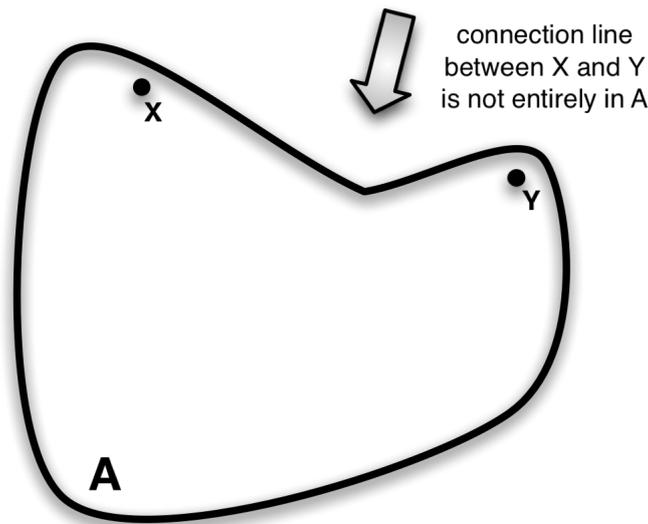


Figure 2.1.: A non-convex set

**Definition 3.** Let  $X$  be convex,  $f : X \rightarrow \overline{\mathbb{R}}$ . We call  $f$  convex if for all  $x, y \in X$  and  $0 < \lambda < 1$ :

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (2.5)$$

We call  $f$  closed if it is lower semi-continuous, i. e., for all  $c \in \mathbb{R}$ ,  $f^{-1}((c, \infty])$  is open in  $X$ . □

Convex sets and functions turn out to be quite convenient:

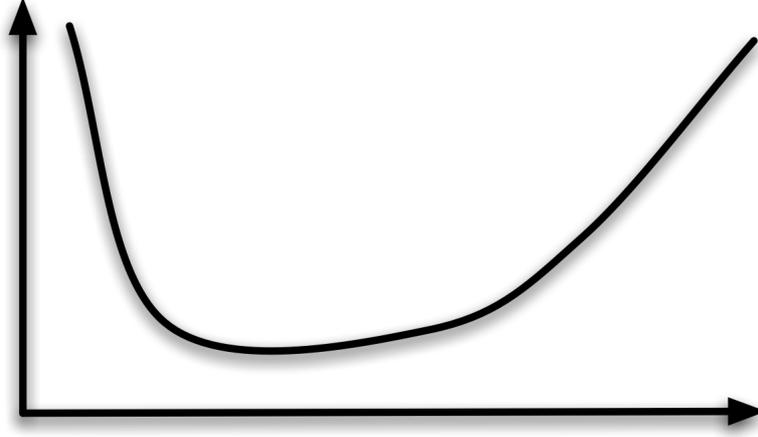


Figure 2.2.: A convex function

**Proposition 1.** *Let  $I$  be an arbitrary set and let  $S_i$  be a convex set for  $i \in I$ . Then the intersection*

$$S = \bigcap_{i \in I} S_i \quad (2.6)$$

*is again convex.*

*Proof.* For all  $x, y \in S$  and all  $i \in I$  holds  $x, y \in S_i$ . □

**Proposition 2.** *Let  $I$  be an arbitrary set. Let  $f_i : X \rightarrow \overline{\mathbb{R}}$  be a convex function and  $\lambda_i \geq 0$  for every  $i \in I$  with at most finitely many  $\lambda_i \neq 0$ . Then holds:*

**(1)** *The function  $f : X \rightarrow \overline{\mathbb{R}}$ ,*

$$f(x) = \sum_{i \in I} \lambda_i f_i(x) \quad (2.7)$$

*is convex. If all  $f_i$  are closed, then so is  $f$ .*

**(2)** *The function  $g : X \rightarrow \overline{\mathbb{R}}$ ,*

$$g(x) = \sup_{i \in I} f_i(x) \quad (2.8)$$

*is convex. If all  $f_i$  are closed, then so is  $g$ .*

**(3)** *If  $(I, \preceq)$  is a directed set and the pointwise limit (proper or improper towards  $+\infty$ ) of  $(f_i)_{i \in I}$  exists, then the function  $h : X \rightarrow \overline{\mathbb{R}}$ ,*

$$h(x) = \lim f_i(x) \quad (2.9)$$

*is convex.*

*Proof.* (1): It is obvious that  $f$  is convex. We show the closedness in two steps ( $\lambda > 0$ ,  $f$  closed  $\Rightarrow \lambda f$  closed;  $f, g$  closed  $\Rightarrow f + g$  closed).

It is  $(\lambda f)^{-1}((c, \infty]) = f^{-1}((c/\lambda, \infty])$ . Let  $x \in (f + g)^{-1}((c, \infty])$ , i. e.,

$$f(x) + g(x) > c. \quad (2.10)$$

We can find an  $\varepsilon > 0$  such that

$$f(x) - \varepsilon + g(x) - \varepsilon > c \quad (2.11)$$

still holds. It is

$$x \in f^{-1}((f(x) - \varepsilon, \infty]) \cap g^{-1}((g(x) - \varepsilon, \infty]) \subset (f + g)^{-1}((c, \infty]) \quad (2.12)$$

with the inner set being open as the intersection of two open sets.

(2): An easy reformulation of the definition of convexity (2.5) is the following: The set

$$E(f) = \{(x, y) \in X \times \mathbb{R} \mid f(x) \leq y\} \quad (2.13)$$

is convex in  $X \times \mathbb{R}$ . ( $E(f)$  is the so-called *epigraph* of  $f$ .) Now, all  $E(f_i)$  are convex and thus,

$$E(g) = \bigcap_{i \in I} E(f_i) \quad (2.14)$$

is also convex by Proposition 1, i. e.,  $g$  is convex. It is  $g$  closed, because

$$g^{-1}((c, \infty]) = \bigcup_{i \in I} f_i^{-1}((c, \infty]) \quad (2.15)$$

is open.

(3): Let  $x, y \in X$ ,  $\lambda \in (0, 1)$ . It is

$$f_i(\lambda x + (1 - \lambda)y) \leq \lambda f_i(x) + (1 - \lambda)f_i(y) \quad (2.16)$$

for all  $i \in I$  and thus, by taking the limit (non-strict inequalities are preserved under limit-taking and the field operations in  $\mathbb{R}$  are continuous),

$$h(\lambda x + (1 - \lambda)y) \leq \lambda h(x) + (1 - \lambda)h(y). \quad (2.17)$$

□

We remark that the limit function  $h$  is *not* closed in general if the  $f_i$  are. This is because even limits of sequences of *continuous* functions can behave nastier than being “just semi-continuous”.

(1) and (3) together imply that also the limit of the series

$$\sum_{i \in I} \lambda_i f_i \quad (2.18)$$

is convex if the sum exists. We can drop the requirement that only finitely many  $\lambda_i$  are nonzero here.

### 2.1.2. Global/Local Minima

Our search for solutions of the minimization problem will become a lot easier with the following theorem. Namely, we can restrict ourselves to search for *local* minima, because all local minima of convex functions are already global.

**Theorem 1.** *Let  $f : X \rightarrow \overline{\mathbb{R}}$  be a convex function and  $x_0 \in X$  with  $f(x_0) \in \mathbb{R}$  which is a local minimum of  $f$ , i. e., there exists a neighbourhood  $U$  of  $x_0$  such that*

$$f(x_0) \leq f(x) \text{ for all } x \in U. \quad (2.19)$$

*Then  $x_0$  is a global minimum of  $f$ .*

*Proof.* Let  $x_1$  be a point in  $X$  with  $f(x_0) > f(x_1)$  and let  $\varepsilon > 0$  such that  $B_\varepsilon(x_0) \subset U$ . It is  $x_1 \notin B_\varepsilon(x_0)$ . We set

$$\lambda = \frac{\varepsilon}{2\|x_1 - x_0\|}. \quad (2.20)$$

Then we have that  $\lambda x_1 + (1 - \lambda)x_0 \in U$  and thus

$$f(x_0) \leq f(\lambda x_1 + (1 - \lambda)x_0) \leq \lambda f(x_1) + (1 - \lambda)f(x_0) < f(x_0) \quad (2.21)$$

which is a contradiction. □

The set of points that minimize a convex function is always convex and closed if the set on which  $f$  is real-valued is closed, i. e.,  $f$  is closed. (The convexity is immediate. Let  $x_0$  be a minimum of  $f$ . The complement of the set of minima is equal to  $\{x \mid f(x) > f(x_0)\}$  which is open, because  $f$  is continuous in the set  $\{x \mid f(x) < \infty\}$ .)

### 2.1.3. Differentiability

From now on, we restrict ourselves to finite-valued functions, i. e., to functions  $f : X \rightarrow \mathbb{R}$ .

Convex functions are always one-sided differentiable in every direction. Because this result is non-trivial, we formulate it as a lemma before using the existence of the differential in the next definition.

**Lemma 1.** Let  $f : X \rightarrow \mathbb{R}$  be convex,  $x \in X^\circ$ , where  $X^\circ$  denotes the interior of  $X$ , and  $v \in H$ . Then the limit

$$\lim_{h \searrow 0} \frac{f(x + hv) - f(x)}{h} \quad (2.22)$$

exists in  $\mathbb{R}$ .

*Proof.* Can be found in [HUL93] p. 238 (for the case  $X = \mathbb{R}^n$ , but the proof for the general case is the same).  $\square$

**Definition 4.** Let  $f : X \rightarrow \mathbb{R}$  be convex,  $x \in X^\circ$  and  $v \in H$ . The real number

$$D_v f(x) = \lim_{h \searrow 0} \frac{f(x + hv) - f(x)}{h} \quad (2.23)$$

is called the directional derivative of  $f$  at  $x$  in direction  $v$ .  $\square$

For a real convex function  $f$ , we essentially have only two directions:  $v = +1$  and  $v = -1$ . The directional derivative  $D_{+1}f(x)$  is equal to the right-sided derivative  $f'_+(x)$  of  $f$  at  $x$ . Likewise,  $D_{-1}f(x)$  is the negative of the left-sided derivative  $-f'_-(x)$ . For arbitrary  $v > 0$ , we have

$$D_v f(x) = \lim_{h \searrow 0} \frac{f(x + hv) - f(x)}{h} \quad (2.24)$$

$$= v \lim_{hv \searrow 0} \frac{f(x + hv) - f(x)}{hv} \quad (2.25)$$

$$= v f'_+(x). \quad (2.26)$$

Analogously, for  $v < 0$ :

$$D_v f(x) = v f'_-(x) \quad (2.27)$$

Of course, we always have ( $v = 0$ ):

$$D_0 f(x) = 0 \quad (2.28)$$

A real function  $g$  is differentiable at  $x$  if and only if  $g'_+(x)$ ,  $g'_-(x)$  exist and are equal. Hence, if  $f$  is differentiable, then

$$D_v f(x) = v f'(x). \quad (2.29)$$

**Definition 5.** Let  $f : X \rightarrow \mathbb{R}$  be convex and  $x \in X^\circ$ . The set

$$\partial f(x) = \{s \in H \mid \langle s, v \rangle \leq D_v f(x) \text{ for all } v \in H\} \quad (2.30)$$

is called the subdifferential of  $f$  at  $x$ .  $\square$

We see immediately that for differentiable real functions  $f$ , we have  $\partial f(x) = \{f'(x)\}$ , because setting  $v$  to  $+1$  resp.  $-1$  in (2.29) and the definition of the subderivative gives us for  $s \in \partial f(x)$

$$s \leq f'(x) \tag{2.31}$$

resp.

$$-s \leq -f'(x), \tag{2.32}$$

hence  $s = f'(x)$ . Moreover,  $f'(x)$  is in  $\partial f(x)$  because of (2.29). This is part of a more general principle:

**Proposition 3.** *Let  $f : X \rightarrow \mathbb{R}$  be convex and  $x \in X$ . Then holds:*

- (1)  $\partial f(x)$  is a nonempty convex compact set.
- (2)  $\partial f(x) = \{s_0\}$  for an  $s_0 \in H$  if and only if  $f$  is differentiable in  $x$ . In this case, we have  $\nabla f(x) = s_0$ .

*Proof.* [HUL93] p. 239. □

Theoretically very interesting, though of no great importance to our problem, is the following

**Theorem 2.** *Let  $X \subset \mathbb{R}^n$  be open and  $f : X \rightarrow \mathbb{R}$  be a convex function. Then  $f$  is differentiable Lebesgue almost everywhere.*

*Proof.* In fact, for  $n = 1$ , it is  $\partial f(x) = [f'_-(x), f'_+(x)]$ . Thus, if  $f$  is not differentiable at  $x$ , then  $\partial f(x)$  is a nontrivial interval and hence contains a rational number in its interior. Further, for  $x < y$ , the interiors of  $\partial f(x)$  and  $\partial f(y)$  are disjoint. This shows that the set of points where  $f$  fails to be differentiable is at most countable and thus Lebesgue zero. For details, refer to [HUL93] pp. 189-190. □

With the notion of the subdifferential, we can formulate our first minimality criterion (which is still very abstract by now).

### 2.1.4. A Minimality Condition

**Theorem 3.** *Let  $f : X \rightarrow \mathbb{R}$  be convex and  $x_0 \in X$ . The following are equivalent.*

- (1)  $x_0$  minimizes  $f$ , i. e.,  $f(x_0) \leq f(x)$  for all  $x \in X$

**(2)**  $0 \in \partial f(x_0)$

*Proof.* ( $\Rightarrow$ ): It is  $f(x_0 + hv) - f(x_0) \geq 0$  for all  $hv$  and thus  $D_v f(x) \geq 0$  for all  $v$ .

( $\Leftarrow$ ): We show that the mapping

$$h \mapsto \frac{f(x_0 + hv) - f(x_0)}{h} \quad (2.33)$$

is nondecreasing for  $h > 0$ . Let  $0 < h_1 < h_2$ . It is

$$f(x_0 + h_1 v) = f\left(\left(1 - \frac{h_1}{h_2}\right)x_0 + \frac{h_1}{h_2}(x_0 + h_2 v)\right) \quad (2.34)$$

$$\leq \left(1 - \frac{h_1}{h_2}\right)f(x_0) + \frac{h_1}{h_2}f(x_0 + h_2 v) \quad (2.35)$$

$$= f(x_0) + h_1 \frac{f(x_0 + h_2 v) - f(x_0)}{h_2} \quad (2.36)$$

which yields

$$\frac{f(x_0 + h_1 v) - f(x_0)}{h_1} \leq \frac{f(x_0 + h_2 v) - f(x_0)}{h_2}. \quad (2.37)$$

This implies that

$$\lim_{h \searrow 0} \frac{f(x_0 + hv) - f(x_0)}{h} = \inf_{h > 0} \frac{f(x_0 + hv) - f(x_0)}{h} \quad (2.38)$$

where the left-hand side is  $\geq 0$  for all  $v \in X$ . This already shows the minimality of  $f(x_0)$ : We can set  $v = x - x_0$ ,  $h = 1$  and get  $f(x) - f(x_0) \geq 0$  with equation (2.38).  $\square$

## 2.2. Bayesian Classification

In supervised learning, the situation is the following. A number of objects (commonly described as real-valued vector) together with a label (a real number or an element of  $\{0, 1\}$ ) are presented to the machine. This set of pairs (object, label) is called the *training set*  $T$ . The machine's objective is to derive a *decision function*  $f$  that maps objects to labels such that it is consistent with the training set (i. e.,  $f(\text{object}) = \text{label}$  for all  $(\text{object}, \text{label}) \in T$ ) and that it predicts labels for new objects *well enough*. In fact, the requirement that  $f$  is consistent with *all* elements of  $T$  is often dropped and replaced by the requirement that it is *close enough* for *most* elements of  $T$ .

For the procedure of finding the decision function  $f$ , two cases can be distinguished ([Vap98]):

1. It is known that the “real” labeling function is taken from a fixed set  $\Gamma = \{f_\alpha \mid \alpha \in A\}$  of functions.
2. No such set is known.

We will be concerned with the first case, i. e., we are given such a set  $\Gamma$  and we only adjust the parameter  $\alpha$  to find the function that fits best. Both approaches assume that such a “real” labeling function exists. One can of course generalize this and search for an appropriate *probability distribution* for the labeling process.

We suppose that the labeling process by the *supervisor* that labeled the training set is determined by a fixed (but unknown) probability distribution  $F$ . Let  $F(\omega|x)$  denote the probability that the supervisor assigns label  $\omega$  to the object  $x$  and let  $F(x)$  denote the probability that object  $x$  is chosen for classification. The problem of finding the best parameter  $\alpha$  is then minimizing the function

$$R(\alpha) = \int L(f_\alpha(x), \omega) dF(x, \omega) \quad (2.39)$$

where  $F(x, \omega) = F(x)F(\omega|x)$  and  $L$  is an appropriately chosen *loss function*, i. e., a non-negative function that increases with the mislabelings of the classification function  $f_\alpha$ . A very simple loss function would be  $L(a, b) = 1 - \delta_{a,b}$  where  $\delta$  denotes the Kronecker symbol, i. e.,  $\delta_{a,b} = 1$  if and only if  $a = b$  and  $= 0$  else.

If we only take in account the only thing we know about the distribution  $F$  – namely the training set  $T = \{(x_1, \omega_1), \dots, (x_\ell, \omega_\ell)\}$  – then (2.39) becomes

$$R'(\alpha) = \frac{1}{\ell} \sum_{i=1}^{\ell} L(f_\alpha(x_i), \omega_i). \quad (2.40)$$

Here we assume that the training set was chosen with the same probability distribution  $F$  as the forthcoming samples.

Since what we do is *binary* classification, i. e., there are only two possible labels for all objects, we can have  $\omega \in \{-1, 1\}$  and by using a simple loss function, we get

$$R''(\alpha) = \sum_{i=1}^{\ell} |f_\alpha(x_i) - \omega_i| \quad (2.41)$$

as the risk function we are to minimize. In fact, up to a constant factor and constant summand, the simple loss function  $(1 - \delta)$  is the only one in binary classification if we demand  $L$  to be symmetric.

## 2.3. Constrained Quadratic Optimization for SVMs

In this and the following section we discuss methods for solving the general optimization problem from definition 1 and explore the problem we are to solve in Support Vector Training. Up to now, we neglected the feasibility set  $S$ , that is we assumed  $S = \mathbb{R}^n$ . As in ordinary vector analysis, finding extremal points in non-open sets is a lot harder than it is in open sets where we have a convenient necessary condition ( $\nabla f = 0$ ). For the case that the feasibility set is a differentiable manifold there is the well-known method of *Lagrangian multipliers* ( $\nabla f = \sum \lambda_i \nabla \varphi_i$ ) for which we will derive a passable substitution (in reality, even a generalization) in the next section. In our case,  $S$  will be a closed convex set and will be described by equality and inequality *constraints*. The feasibility set will be expressed as the intersection of inverse images of closed sets under convex functions.

Suppose our training set is

$$T = \{(x_1, \omega_1), \dots, (x_\ell, \omega_\ell)\} \quad (2.42)$$

where  $\omega_i \in \{-1, +1\}$  and  $x_i \in \mathbb{R}^n$  for  $1 \leq i \leq \ell$ . We split  $T$  into the two sets

$$A = \{x_i \mid (x_i, \omega_i) \in T, \omega_i = +1\}, \quad (2.43)$$

$$B = \{x_i \mid (x_i, \omega_i) \in T, \omega_i = -1\}. \quad (2.44)$$

Our goal is to find an affine hyperplane  $V$  in  $\mathbb{R}^n$  such that the points of  $A$  are on one side and the points of  $B$  are on the other. Of course, this is only possible if  $A \cap B = \emptyset$ , i. e., no point has *both* labels. Additionally, this hyperplane should have maximal distance to the points of  $A \cup B$ . More formally, if  $V$  is described by the equation

$$\langle x, w \rangle = c \quad (2.45)$$

where  $w \in \mathbb{R}^n$ ,  $\|w\| = 1$  and  $c \in \mathbb{R}$ , then it should hold that  $\langle x, w \rangle > c$  for  $x \in A$  and  $\langle y, w \rangle < c$  for  $y \in B$ . In this case, we say that  $V$  *separates the sets  $A$  and  $B$*  and that  $A$  and  $B$  are *separable*. The distance  $\text{dist}(x, V)$  of a point  $x$  to  $V$  is equal to  $|c - \langle x, w \rangle|$ . More explicitly for our points of interest,

$$\text{dist}(x, V) = \begin{cases} \langle x, w \rangle - c & , x \in A \\ c - \langle x, w \rangle & , x \in B \end{cases} \quad (2.46)$$

The following exposition closely follows [Vap98], chap. 10 and [Bur98], chap. 3. We define for every normed  $w$  the numbers  $c_1(w) = \inf\{\langle x, w \rangle \mid x \in A\}$  and  $c_2(w) = \sup\{\langle y, w \rangle \mid y \in B\}$ . According to (2.46), the sum of the minimal distances of  $A$  and  $B$  to  $V$  respectively is equal to  $\rho(w) = (c_1(w) - c) + (c -$

$c_2(w)) = c_1(w) - c_2(w)$ . We note that even though the parameter  $c$  is yet to be determined, the expression  $\rho(w)$  does not depend on it. Since we want the distances to be maximal, we want to maximize  $\rho(w)$ . After we found such a normed  $w$  for which  $\rho(w)$  is maximal, it is clear that  $c = (c_1(w) + c_2(w))/2$  is the optimal choice for the parameter  $c$ . For it is the mean of both extreme values.

**Lemma 2.** *The function  $w \mapsto \rho(w)$  has a unique maximum if  $A, B \neq \emptyset$  and  $A \cap B = \emptyset$ .*

*Proof.* The existence of a maximum is clear since  $S^{n-1}$  is compact and  $\rho$  is continuous. To prove the uniqueness, we first show that  $\rho$ , as a function on the set of  $x \in \mathbb{R}^n$  for which  $\|x\| \leq 1$ , attains its maximum at the set's boundary, i. e.,  $S^{n-1}$ . For let  $0 < \|x\| < 1$ , then  $\rho(x/\|x\|) = \rho(x)/\|x\| > \rho(x)$  with  $\|x/\|x\|\| = 1$ . Let now  $x_1, x_2 \in S^{n-1}$  be two distinct maxima. Then  $(x_1 + x_2)/2$  is also a maximum with  $\|(x_1 + x_2)/2\| < 1$ . Contradiction.  $\square$

We can reformulate the problem to finding  $w \in \mathbb{R}^n \setminus \{0\}$  and  $b \in \mathbb{R}$  such that

$$\langle x, w \rangle - b \geq +1 \quad (x \in A) \tag{2.47}$$

$$\langle y, w \rangle - b \leq -1 \quad (y \in B) \tag{2.48}$$

where  $\|w\|$  is minimal. The above is equivalent to

$$\omega_i (\langle x_i, w \rangle - b) - 1 \geq 0 \quad (1 \leq i \leq \ell). \tag{2.49}$$

If we find an optimal  $w$ , equality holds in (2.49) for at least two  $i$  – one of each class. Then, the distance between the two hyperplanes defined by  $\langle x, w \rangle - b = \pm 1$  is equal to  $2/\|w\|$ . These hyperplanes are parallel to  $V$  which is defined by  $\langle x, w \rangle - b = 0$  and are called *support hyperplanes*. Thus, the above found value  $2/\|w\|$  is the same as that of the function  $\rho$  for the respective argument. To sum up, our problem to solve is now the following.

$$\min \|w\|^2 \quad \text{subject to (2.49)} \tag{2.50}$$

This problem is a convex optimization problem since the objective function  $x \mapsto \|x\|^2$  is convex and so is the feasibility set  $S$  which is the intersection of half spaces defined by (2.49).

## 2.4. Are we there yet? – Optimality Conditions

In this section, we will introduce necessary and sufficient conditions for minima of convex optimization problems. These are known as the *Karush-Kuhn-Tucker (KKT)* conditions.

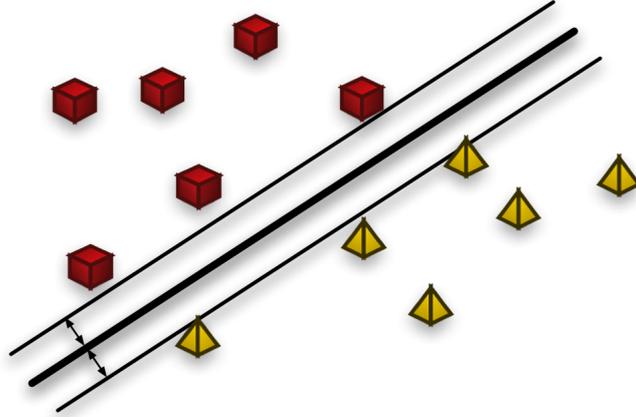


Figure 2.3.: Support hyperplanes

**Theorem 4 (KKT).** Let  $f_i : X \rightarrow \mathbb{R}$  be convex functions for  $0 \leq i \leq m$  where  $X \subset \mathbb{R}^n$  is a convex set. We define

$$S = \{x \in X \mid f_i(x) \leq 0 \quad (1 \leq i \leq m)\}. \quad (2.51)$$

Suppose that  $x^*$  minimizes  $f_0$  in the set  $S$ . Then there exists  $(\lambda_0^*, \lambda^*) = (\lambda_0^*, \dots, \lambda_m^*) \neq 0$  such that

- (1) the function  $L(x, \lambda_0^*, \lambda^*) = \sum_{i=0}^m \lambda_i^* f_i(x)$  is minimized by  $x^*$  in the set  $X$ .
- (2)  $(\lambda_0^*, \lambda^*) \geq 0$ .
- (3)  $\lambda_i^* f_i(x^*) = 0$  for all  $i$ .

Let now  $x^*$  and  $(\lambda_0^*, \lambda^*)$  satisfy conditions (1), (2), (3). If  $\lambda_0^* \neq 0$ , then  $x^*$  minimizes  $f_0$  in  $S$ . □

The function

$$L(x, \lambda_0, \lambda) = \sum_{i=0}^m \lambda_i f_i(x) \quad (2.52)$$

is called the *Lagrangian* of the optimization problem.

**Proposition 4.** With the notation of theorem 4, it is sufficient for  $\lambda_0^*$  to be possibly  $\neq 0$  that there exists an  $x_0 \in X$  with  $f_i(x_0) < 0$  for all  $1 \leq i \leq m$ . □

Thus, in this case, we can eliminate  $\lambda_0$  from our Lagrangian by division and get the simpler version

$$L(x, \lambda) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x). \quad (2.53)$$

## 2.5. Kernel Functions

It is possible to replace the separating hyperplane by a more general hypersurface in  $\mathbb{R}^n$  which is the inverse image of a linear hyperplane in a higher (often infinite) dimensional Hilbert space  $H$ . The technique we will use is known as the *kernel trick*. It enables us to classify data sets that are not separable by a (linear) hyperplane.

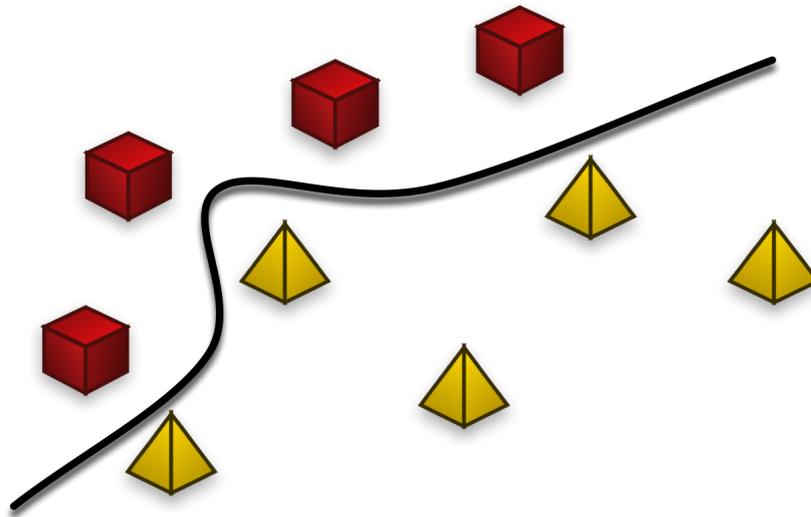


Figure 2.4.: Data that is not linearly separable

We consider a transformation  $\Phi : \mathbb{R}^n \rightarrow H$  with which we map the training data before applying the algorithm. Thus, the inner products that occur become  $\langle \Phi(x), \Phi(y) \rangle = K(x, y)$ . We see that we do not really have to know the mapping  $\Phi$  – not even the space  $H$ . We only need to know the *kernel function*  $K$ .

Now, one can go the opposite direction and ask: For which functions  $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  exists a Hilbert space  $(H, \langle \cdot, \cdot \rangle_H)$  and a transformation  $\Phi : \mathbb{R}^n \rightarrow H$ , such that  $K(x, y) = \langle \Phi(x), \Phi(y) \rangle_H$ ?

This question is partially answered by *Mercer's theorem* (in [Mer09], although this version is a generalization of his original theorem to more general domains than the real compact intervals  $[a, b]$ ). It states that this is true for all continuous symmetric  $K$  that are nonnegative definite. The following proof can be skipped.

**Theorem 5** (Mercer). *Let  $(X, \mathfrak{A}, \mu)$  be a  $\sigma$ -finite measure space (i. e., for  $\iota \in \mathbb{N}$  there exist  $A_\iota \in \mathfrak{A}$  with  $\mu(A_\iota) < \infty$  such that  $X = \bigcup A_\iota$ ). Let further  $K :$*

$X \times X \rightarrow \mathbb{R}$  be a symmetric function in  $L^2(X^2)$  such that for all  $f \in L^2(X)$  there holds:

$$\int_{X^2} K(x, y) f(x) f(y) d(x, y) \geq 0 \quad (2.54)$$

Then there exists an orthonormal family  $(f_i)_{i \in I}$  in  $L^2(X)$  and a family  $(\lambda_i)_{i \in I}$  of non-negative real numbers such that

$$K(x, y) = \sum_{i \in I} \lambda_i f_i(x) f_i(y) \quad (2.55)$$

for almost all  $(x, y) \in X^2$ .

*Proof.* We define the operator  $T = T_K : L^2(X) \rightarrow L^2(X)$  by

$$T_K f(x) = \int_X K(x, t) f(t) dt. \quad (2.56)$$

It maps into  $L^2(X)$ , because  $K$  is in  $L^2(X^2)$ . For  $L^2(X)$  is a Hilbert space, there exists an orthonormal basis  $B$  of  $L^2(X)$ . We will show that  $\sum \|Tb\|^2$  is finite if  $b$  varies in  $B$ :

$$\sum_{b \in B} \|Tb\|^2 = \sum_{b \in B} \int_X \left| \int_X K(x, t) b(t) dt \right|^2 dx \quad (2.57)$$

$$= \int_X \sum_{b \in B} |\langle K(x, \cdot), b \rangle|^2 dx \quad (2.58)$$

$$= \int_X \|K(x, \cdot)\|^2 dx \quad (2.59)$$

$$= \int_X \int_X |K(x, y)|^2 dy dx < \infty \quad (2.60)$$

This shows that  $T$  is Hilbert-Schmidt and hence compact ([Wei00], Satz 3.18(b)). Since  $K$  is symmetric, so is  $T$ . We can thus apply the spectral theorem and get the existence of an orthonormal basis  $(f_i)_{i \in I}$  of  $L^2(X)$  which consists of eigenfunctions of  $T$ . Let  $Tf_i = \lambda_i f_i$  for  $i \in I$ . It is  $\lambda_i = \lambda_i \langle f_i, f_i \rangle = \langle Tf_i, f_i \rangle \geq 0$  for all  $i$ . Further, for  $f \in L^2(X)$ , it is  $\int f(t) \sum \lambda_i f_i(t) f_i(x) dt = \sum \lambda_i f_i(x) \int f_i(t) f(t) dt = \sum f_i(x) \langle T_K f_i, f \rangle = \sum \langle T_K f, f_i \rangle f_i(x) = T_K f(x)$  almost everywhere. Now, the mapping  $K \mapsto T_K$  is injective since  $X$  is  $\sigma$ -finite. Hence, the claimed formula follows.  $\square$

We can then define the transformation  $\Phi : X \rightarrow \ell^2(I)$  by

$$\Phi(x) = \left( \sqrt{\lambda_i} f_i(x) \right)_{i \in I} \quad (2.61)$$

where our Hilbert space  $H = \ell^2(I)$  is the space of quadratic summable real sequences with index set  $I$  equipped with the inner product

$$\langle a, b \rangle = \sum_{i \in I} a_i b_i \quad (2.62)$$

where  $a = (a_i)$ ,  $b = (b_i)$ . The image  $\Phi(x)$  is really in  $\ell^2(I)$ , because

$$\sum_{i \in I} \left| \sqrt{\lambda_i} f_i(x) \right|^2 = \sum_{i \in I} \lambda_i f_i(x) f_i(x) = K(x, x). \quad (2.63)$$

In particular,  $\Phi(x)$  is quadratic summable. This gives us our desired result

$$K(x, y) = \langle \Phi(x), \Phi(y) \rangle. \quad (2.64)$$

Our favored space  $\mathbb{R}^n$  is  $\sigma$ -finite (with respect to the Lebesgue measure  $\lambda$ ). We can choose  $A_\iota = \{x \in \mathbb{R}^n \mid \|x\| < \iota\}$ . It is  $\lambda(A_\iota) = \iota^n \pi^{n/2} / \Gamma(\frac{n}{2} + 1) < \infty$  and  $A_\iota \rightarrow \mathbb{R}^n$ .

Possible nonnegative definite kernels include ([Bur98])

- $K(x, y) = (\langle x, y \rangle + 1)^p \quad (p \in \mathbb{N})$
- $K(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}} \quad (\sigma \neq 0)$
- $K(x, y) = \tanh(\kappa \langle x, y \rangle - \delta)$  for *certain*  $\kappa, \delta \in \mathbb{R}$

where we might have to restrict  $K$  to a smaller set than the whole  $\mathbb{R}^n$ . For example, the *Gaussian* kernel  $K(x, y) = e^{-\|x-y\|^2/2\sigma^2}$  is not in  $L^2(\mathbb{R}^n \times \mathbb{R}^n)$ :

$$\int_{\mathbb{R}^n \times \mathbb{R}^n} |K(x, y)|^2 d(x, y) = \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} e^{-\|x-y\|^2/\sigma^2} dx dy \quad (2.65)$$

$$= \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} e^{-\|x\|^2} |\sigma|^n dx dy \quad (2.66)$$

$$= \int_{\mathbb{R}^n} |\sigma \sqrt{\pi}|^n dy = \infty \quad (2.67)$$

But it is, of course, in  $L^2(C \times C)$  for every compact subset  $C$  of  $\mathbb{R}^n$ . ( $\int K^2 \leq \lambda(C)^2 \max K^2$ )

## 3. SVM Algorithms

This chapter introduces the usual formulation of the SVM training problem and reviews popular solution algorithms. Caching and shrinking techniques are also treated though they are not applicable for use on the microcontroller due to the stringent memory space limitations. In fact, one could skip all the preceding and start with section 3.3, the final optimization problem statement, if one is only interested in the algorithmic aspects of Support Vector Machines as opposed to the mathematical aspects and derivations.

### 3.1. Naïve KKT

With our convex constraint functions (2.49),  $f_i(w, b) = -\omega_i(\langle x_i, w \rangle - b) + 1$  for  $1 \leq i \leq \ell$ , and a slightly modified objective function  $f_0(w, b) = \frac{1}{2} \|w\|^2$ , the simple Lagrangian (2.53) becomes

$$L(w, b, \lambda) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^{\ell} \lambda_i \omega_i (\langle x_i, w \rangle - b) + \sum_{i=1}^{\ell} \lambda_i. \quad (3.1)$$

The KKT conditions tell us that it is necessary for  $(w, b)$  to be a solution of (2.50) that  $(w, b)$  is a minimum of the Lagrangian for a certain choice of  $\lambda$ . Since  $L$  is continuously differentiable, it is therefore necessary that the partial derivatives  $\partial L / \partial w$  and  $\partial L / \partial b$  vanish. This means that  $w - \sum \lambda_i \omega_i x_i = 0$ , which is

$$w = \sum_{i=1}^{\ell} \lambda_i \omega_i x_i. \quad (3.2)$$

Also, by partial derivation with respect to  $b$ ,

$$\sum_{i=1}^{\ell} \lambda_i \omega_i = 0. \quad (3.3)$$

We can substitute this into the Lagrangian by noticing

$$\frac{1}{2} \|w\|^2 = \frac{1}{2} \left\langle \sum \lambda_i \omega_i x_i, \sum \lambda_i \omega_i x_i \right\rangle = \frac{1}{2} \sum_{i,j=1}^{\ell} \lambda_i \lambda_j \omega_i \omega_j \langle x_i, x_j \rangle \quad (3.4)$$

and

$$\sum_{i=1}^{\ell} \lambda_i \omega_i (\langle x_i, w \rangle + b) = \sum_{i,j=1}^{\ell} \lambda_i \lambda_j \omega_i \omega_j \langle x_i, x_j \rangle + 0 \quad (3.5)$$

which yields

$$W(\lambda) = L(w, b, \lambda) = \sum_{i=1}^{\ell} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \lambda_i \lambda_j \omega_i \omega_j \langle x_i, x_j \rangle. \quad (3.6)$$

This new formulation of the Lagrangian does not depend on  $w$  and  $b$  anymore. Since it is necessary that  $(w, b)$  minimizes  $L$  for a choice of  $\lambda$  subject to certain constraints and since we know that a minimum exists, it is sufficient to maximize  $W$  with respect to  $\lambda$  subject to the same constraints. This reformulation, the so-called *dual* formulation, is summarized in section 3.3.

## 3.2. Dealing with Errors

There may be cases where we want to tolerate some training errors, i. e., points of the training set that lie on the wrong side of the hyperplane. This is the case, for example, if we use a linear classifier and the subsets of the training set that correspond to the respective labels are not linearly separable (i. e., there exists no hyperplane that separates the two sets). To achieve this, we introduce a penalty parameter  $C$  for points that fail to be on the side of its label. Of course, we want the penalty to be higher the greater the distance of the erroneous points to the hyperplane. But first of all, we need to loosen the strict constraint  $\omega_i (\langle x_i, w \rangle + b) \geq 1$ . We therefore introduce non-negative slack variables  $\xi_i$  ([Bur98], 3.5) to transform the above into

$$\omega_i (\langle x_i, w \rangle - b) \geq 1 - \xi_i \quad (1 \leq i \leq \ell) \quad (3.7)$$

where we want to have  $\xi \geq 0$ . To actually implement the penalty, we simply add the slack variables to the objective function of the minimization problem multiplied with the parameter  $C$ :

$$f_0(w, b, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \quad (3.8)$$

What changes does this introduce into the dual formulation? Well, the full Lagrangian reads as follows (note the changed/additional constraints in the primal formulation!).

$$L(w, b, \xi, \lambda, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i - \sum_{i=1}^{\ell} \lambda_i (\omega_i (\langle x_i, w \rangle - b) - 1 + \xi_i) - \sum_{i=1}^{\ell} \mu_i \xi_i \quad (3.9)$$

Again, this can be simplified by setting  $\partial L/\partial w = \partial L/\partial b = 0$ . These equations yield the same as above. Additionally, we can have  $\partial L/\partial \xi_i = 0$  for  $1 \leq i \leq \ell$  which is  $C - \lambda_i - \mu_i = 0$ . This implies

$$C \sum_{i=1}^{\ell} \xi_i - \sum_{i=1}^{\ell} \lambda_i \xi_i - \sum_{i=1}^{\ell} \mu_i \xi_i = 0. \quad (3.10)$$

Thus, in reality  $\xi$  and  $\mu$  do not appear at all in the dual Lagrangian  $W(\lambda)$ . The only additional constraint we get, since  $\mu_i \geq 0$ , is  $C \geq \lambda_i$ .

### 3.3. The Dual Formulation

We already utilize section 2.5 here, i. e., we replace the scalar product  $\langle x, y \rangle$  by a kernel function  $K(x, y)$ . Further, we again state it as a minimization problem by flipping signs of the objective function. The complete dual formulation then reads:

$$\min \quad W(\lambda) = - \sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i,j=1}^{\ell} \lambda_i \lambda_j \omega_i \omega_j K(x_i, x_j) \quad (3.11)$$

$$\text{s. t.} \quad \sum_{i=1}^{\ell} \lambda_i \omega_i = 0 \quad (3.12)$$

$$0 \leq \lambda \leq C \quad (3.13)$$

In the above formula,  $0 \leq \lambda \leq C$  means  $0 \leq \lambda_i \leq C$  for all  $i$ . With the notation  $Q = (\omega_i \omega_j K(x_i, x_j))_{1 \leq i, j \leq \ell}$  and  $e = (1)_{1 \leq i \leq \ell}$ , it becomes

$$\min \quad W(\lambda) = -\lambda^T e + \frac{1}{2} \lambda^T Q \lambda \quad (3.14)$$

$$\text{s. t.} \quad \lambda^T \omega = 0 \quad (3.15)$$

$$0 \leq \lambda \leq C e \quad (3.16)$$

where  $z^T$  denotes the transpose of  $z$ . The necessary and sufficient KKT conditions for a minimum are ([Pla98]):

$$\lambda_i = 0 \iff \omega_i u_i \geq 1 \quad (3.17)$$

$$0 < \lambda_i < C \iff \omega_i u_i = 1 \quad (3.18)$$

$$\lambda_i = C \iff \omega_i u_i \leq 1 \quad (3.19)$$

Here  $u_i$  denotes the “raw” classifier function evaluated at the training point  $x_i$ , that is

$$u_i = \sum_{j=1}^{\ell} \lambda_j \omega_j K(x_i, x_j) - b. \quad (3.20)$$

### 3.4. Solution Strategies

In this section, we briefly discuss popular algorithms for SVM training. Note that these are geared towards very large sets of input data and thus mostly not immediately utilizable for microcontroller use where we do not expect such magnitude of data. Also, these algorithms assume that almost an infinite amount of data can temporarily be stored on a hard disk drive and that only RAM space is limited, which is not true in an MCU environment.

#### 3.4.1. Osuna et al.

This algorithm ([OFG97]) utilizes a *decomposition* of the input index set  $\{1, \dots, \ell\}$  into the *working set*  $B$  and the remainder set  $N$ , whose associated multipliers  $\lambda_i$  will not change in the current iteration. If we denote by  $\lambda_J$ ,  $\omega_J$  and  $Q_{IJ}$  the vectors and matrices with entries corresponding to the index sets  $I, J \subset \{1, \dots, \ell\}$ , then the optimization problem becomes

$$\min \quad -\lambda_B^T e + \frac{1}{2} \lambda_B^T Q_{BB} \lambda_B + \lambda_B^T q_{BN} \quad (3.21)$$

$$\text{w. r. t.} \quad \lambda_B \quad (3.22)$$

$$\text{s. t.} \quad \lambda_B^T \omega_B + \lambda_N^T \omega_N = 0 \quad (3.23)$$

$$0 \leq \lambda_B \leq Ce \quad (3.24)$$

Here, we omitted the constant terms that only include  $\lambda_N$  and  $Q_{NN}$ . Further,  $q_{BN} = \left( \omega_i \sum_{j \in N} \lambda_j \omega_j K(x_i, x_j) \right)_{i \in B}$ . The algorithm is now based on the following two observations.

- If we move an arbitrary index  $i$  from  $B$  to  $N$ , the objective function (of the original problem) does not change and the solution is feasible. (*Build down*)
- If we move an index  $i$  from  $N$  to  $B$  that violates the KKT conditions and solve the subproblem for  $B$ , there is a strict improvement of the objective function. (*Build up*)

The algorithm is sketched in pseudo-code below in figure 3.1.

Despite its good reception and reported positive results, the algorithm has a theoretical disadvantage: Though it is guaranteed that the solution improves in each iteration, there is no proof that it actually converges to an optimal solution ([CHL00]).

---

```

Osuna( $x, \omega$ ) {
  choose  $B \subset \{1, \dots, \ell\}$  arbitrarily;
   $N := \{1, \dots, \ell\} \setminus B$ ;
  for(;;)
  {
    solve subproblem for  $B$ ;
    if( $\exists i \in N$  such that  $\lambda_i$  violates KKT)
    {
      choose  $j \in B$  arbitrarily;
       $B := \{i\} \cup B \setminus \{j\}$ ;
       $N := \{j\} \cup N \setminus \{i\}$ ;
    } else break;
  }
  return  $\lambda$ ; }

```

---

Figure 3.1.: Osuna et. al. decomposition algorithm

### 3.4.2. SMO

Sequential Minimal Optimization (SMO, [Pla98]) essentially employs the idea of Osuna's decomposition algorithm with  $|B| = 2$  and adds heuristics for the choice of the next working set pair. The main advantage of having only two multipliers in the working set at a time is that the optimal solution can be computed analytically here and the algorithm therefore does not have to rely on the usage of numeric quadratic program solvers. We will take more time to explain and derive this method since we will use this algorithm in the implementation (chapter 4).

We consider the two Lagrangian multipliers  $\lambda_1$  and  $\lambda_2$ . (We assume  $B = \{1, 2\}$  without loss of generality.) The constraint (3.16) is  $0 \leq \lambda_1, \lambda_2 \leq C$  and (3.15) means  $\omega_1 \lambda_1 + \omega_2 \lambda_2 = \omega_1 \lambda'_1 + \omega_2 \lambda'_2$  where  $\lambda'_i$  denotes the old value of  $\lambda_i$  of the previous iteration. Following [Pla98], we first compute the optimal value for  $\lambda_2$  and then calculate  $\lambda_1$  from the constraints. We distinguish the cases  $\omega_1 = \omega_2$  and  $\omega_1 \neq \omega_2$ . In the first case, we have  $\lambda_1 + \lambda_2 = d$ , in the second  $\lambda_1 - \lambda_2 = d$  where  $d$  is a constant. Thus, the possible values for  $(\lambda_1, \lambda_2)$  lie all on a line segment depicted in figure 3.2.

The lower and upper limits for  $\lambda_2$  thus are:  $L = \max(0, \lambda'_2 + \lambda'_1 - C)$ ,  $H = \min(C, \lambda'_2 + \lambda'_1)$  for  $\omega_1 = \omega_2$  and  $L = \max(0, \lambda'_2 - \lambda'_1)$ ,  $H = \min(C, C + \lambda'_2 - \lambda'_1)$

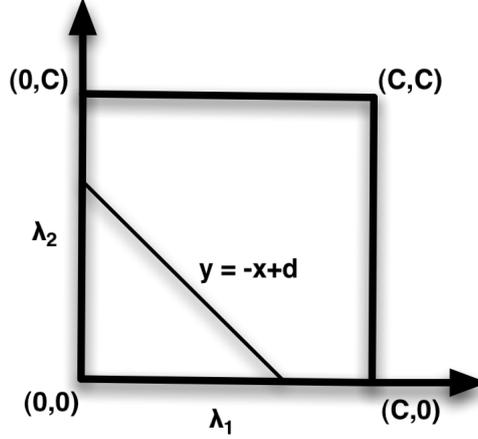


Figure 3.2.: Feasibility line for  $(\lambda_1, \lambda_2)$  in the case  $\omega_1 = \omega_2$

for  $\omega_1 \neq \omega_2$ . The objective function (c. f. Osuna) is

$$-\lambda_1 - \lambda_2 + \frac{1}{2}K_{11}\lambda_1^2 + \frac{1}{2}K_{22}\lambda_2^2 + sK_{12}\lambda_1\lambda_2 + \lambda_1v_1 + \lambda_2v_2 \quad (3.25)$$

where  $K_{ij} = K(x_i, x_j)$ ,  $v_i = \sum_{j=3}^{\ell} \lambda'_j \omega_j K_{ij} = u_i + b' - \lambda'_1 \omega_1 K_{1i} - \lambda'_2 \omega_2 K_{2i}$  and  $s = \omega_1 \omega_2 = \pm 1$  depending on whether  $\omega_1 = \omega_2$ . By using  $\lambda_1 + s\lambda_2 = \lambda'_1 + s\lambda'_2 = d$ , we can transform this into a function that depends on  $\lambda_2$  only. We can then set the ordinary first derivative equal to zero and calculate  $\lambda_2$ , provided  $\eta = K_{11} + K_{22} - 2K_{12}$ , which is the second derivative of this function, does not vanish. The embarrassing situation that it *does* vanish can occur, for example, if  $x_i = x_j$  for  $i \neq j$ . In the other case, the optimal  $\lambda_2$  is equal to

$$\lambda_2^{\text{new}} = \lambda'_2 + \frac{\omega_2(u_1 - \omega_1 - u_2 + \omega_2)}{\eta}. \quad (3.26)$$

The quantity  $E_i = u_i - \omega_i$  is called the error of the  $i$ th training example. Next we need to check whether  $\lambda_2^{\text{new}} \in [L, H]$  and clip it into our square if it lies outside:

$$\lambda_2^{\text{new,clipped}} = \begin{cases} H & , \lambda_2^{\text{new}} > H \\ \lambda_2^{\text{new}} & , L \leq \lambda_2^{\text{new}} \leq H \\ L & , \lambda_2^{\text{new}} < L \end{cases} \quad (3.27)$$

We can then compute  $\lambda_1$  from our equality constraint:

$$\lambda_1 = \lambda'_1 + s \left( \lambda'_2 - \lambda_2^{\text{new,clipped}} \right) \quad (3.28)$$

If  $\eta = 0$ , we just evaluate the objective function at the boundaries  $\lambda_2 = L$  and  $\lambda_2 = H$  and check whether the values differ. If so, we take the lower value. If not, then  $\lambda_2 = \lambda'_2$  and we cannot make any progress here.

Two heuristics are used to determine the working set pair for the next iteration. The first heuristic is concerned with finding a suitable  $\lambda_2$  and utilizes the fact that many multipliers end up being either 0 or  $C$  by termination of the algorithm. Thus, after a first pass through all examples, the algorithm consequently only chooses training examples where the corresponding multiplier is *strictly between* 0 and  $C$ . When there are no more changes possible with these examples, it returns to looping over *all* examples again. The second heuristic chooses a suitable partner for  $\lambda_2$ , i. e., one with the largest expected step size. To approximate the expected step size, the training example errors  $E_i$  are used. These are stored in an *error cache* for fast access. The choice with largest expected progress is the one where  $|E_1 - E_2|$  is maximal. If there is no progress with this example, the algorithm first loops over all examples currently not at the bounds and then over all examples until progress is made.

### 3.4.3. SVM<sup>light</sup>

T. Joachims introduced two new methods for solving the SVM training problem in [Joa98], where he presented the SVM<sup>light</sup> package. The first one is a more sophisticated method for selecting the working set than the “random” one used in Osuna decomposition. He proposed using first-order approximation to the objective function to find a direction  $d$  of steepest descent, in which the algorithm continues its operation. For this, he solves the following problem:

$$\min \quad V(d) = (\nabla W(\lambda))^T d \quad (3.29)$$

$$\text{s. t.} \quad \omega^T d = 0 \quad (3.30)$$

$$d_i \geq 0 \quad (\lambda_i = 0) \quad (3.31)$$

$$d_i \leq 0 \quad (\lambda_i = C) \quad (3.32)$$

$$-e \leq d \leq e \quad (3.33)$$

$$|\{i \mid d_i \neq 0\}| = q \quad (3.34)$$

Here  $q = |B|$  is the size of the working set. Our new working set will then be  $\{i \mid d_i \neq 0\}$ . Joachims gave an easy way to compute the solution of this optimization problem by sorting the  $\lambda_i$  in a clever way. The second method is “shrinking” – a technique that, like SMO’s heuristic, uses the fact that there are many multipliers at the bounds in the optimal solution to reduce the size of the optimization problem. Of course, if the guess that a multiplier will be at the bounds was wrong, then it has to be visited in a later iteration nonetheless.

### 3.5. The Classifier Function

When we are done with the training algorithm and found our optimal  $\lambda$  and  $b$ , we are bound to ask how we can use this knowledge in classification of new examples. In the case of a *linear* SVM, i. e., no kernel function but the ordinary scalar product was used, we can simply calculate the vector  $w$  by

$$w = \sum_{i=1}^{\ell} \lambda_i \omega_i x_i \quad (3.35)$$

and the classification function is

$$f(x) = \text{sgn}(\langle x, w \rangle - b). \quad (3.36)$$

The situation is not that easy if a kernel function was used. We cannot calculate  $w$ , because we do not even know the space  $H$  it belongs to. We therefore have to expand  $w$  in (3.36) to get

$$f(x) = \text{sgn} \left( \sum_{i=1}^{N_s} \lambda_i \omega_i K(x, x_i) - b \right) \quad (3.37)$$

where  $N_s$  denotes the number of *support vectors*, i. e., the number of vectors for which  $\lambda_i \neq 0$ . We assume here without loss of generality that the  $x_i$  are numbered in such a way that the first  $N_s$  vectors are the support vectors. These are the only ones we need to remember after the training process.

## 4. Implementation

This chapter describes the implementation of  $\mu$ SVM, a Support Vector Machine package for use on small microcontroller units.

### 4.1. $\mu$ SVM Overview

Sequential Minimal Optimization (section 3.4.2) is utilized by the  $\mu$ SVM package for solving the quadratic program in SVM training. It supports the data types `char`, `int` and `float` for training example vectors. Support for new data types can easily be implemented with minimal changes. The decision for a specific type is made at compile time by macro definitions (e. g., compiler flag `-DuSVM_X_FLOAT` for `float` vectors, see Documentation for details). Different and also user-added kernel functions can be used for the algorithm. It can be changed at run-time by setting the function pointer `uSVM_ker`. Such a change can be useful, for example, if the timing requirements of the program change at run-time and the training procedure has to terminate earlier than in normal operation. Then, a faster kernel could be used. Also, the precision `uSVM_EPS` can always be changed. Since the final values of the Lagrangian multipliers are approximated fairly good early in the course of the training algorithm (see chapter 5), the value of the flag `uSVM_terminate` is checked periodically. The current multiplier values are output immediately and the algorithm is stopped in the case the flag was set. The main data structures are:

- `uSVM_x`: Pointer to the training example vectors. The dimension of the vectors is given by `uSVM_n` and the number of examples is stored in `uSVM_e11`. This field is accessed by the `uSVM_READ(i,k)` and `uSVM_WRITE(i,k,z)` macros.
- `uSVM_omega`: This array contains the labels of the example vectors. The values here should only be  $\pm 1$ . The length of the array is again `uSVM_e11`.
- `uSVM_lambda`: Array of Lagrangian multipliers with `float` precision.
- `E`: The error cache described together with the SMO algorithm. It is used to determine the next working set pair. It also speeds up the train-

ing process by reducing the number of necessary kernel evaluations dramatically. Since this array consists of `uSVM_e11` floating point values, the memory requirements of the training process nearly doubles when using the error cache. It can therefore be deactivated by setting the `uSVM_NO_ERROR_CACHE` macro.

The relevant functions and procedures are:

- `uSVM_train()`: Starts the training algorithm. Allocates the `uSVM_lambda` field which is *not free()*'d by the function itself, but has to be freed by the application in case it is no longer used. Returns `-1` in case of an error and the number of support vectors otherwise. This number is also stored in the `uSVM_nSV` variable.
- `examine(i2)`: Searches for a suitable working set partner for `i2` until either progress is made or all examples were tried.
- `take_step(i1,i2)`: Computes the optimal solution for the subproblem induced by the indices `i1` and `i2`. Returns `1` if the current solution was improved and `0` otherwise.
- `uSVM_classify(x)`: Classifies the example given by the array `x` of dimension `uSVM_n`. Returns  $\pm 1$ .

The training algorithm is sketched in figure 4.1 on page 29.

## 4.2. Target Hardware

$\mu$ SVM was developed and tested on an Atmel AVR ATmega16 microcontroller using `avr-libc` version 1.2.5. The ATmega16 has a 16 MHz pipelined RISC processor with 1 kB internal RAM. The register size is 8 bit.

---

```

uSVM_train() {
    uSVM_lambda = malloc();
    while progress was made in previous iteration
        for all indices i2
            examine(i2);
    delete non-support vectors from uSVM_x, uSVM_omega and uSVM_lambda;
    nSV = # of support vectors;
    return nSV;
}

examine(i2) {
    while not tried all examples i1
        i1 = good working set partner for i2;
        if (take_step(i1,i2)==1) return 1;
    return 0;
}

take_step(i1,i2) {
    (lambda1,lambda2) = optimal solution for i1 and i2;
    if(lambda1==uSVM_lambda[i1] && lambda2==uSVM_lambda[i2])
        return 0; // no progress
    update threshold b;
    update error cache E;
    uSVM_lambda[i1] = lambda1;
    uSVM_lambda[i2] = lambda2;
    return 1;
}

```

---

Figure 4.1.:  $\mu$ SVM training algorithm

## 5. Results and Discussion

This chapter evaluates the temporal behavior and the numerical accuracy of the  $\mu$ SVM package. Also, it is examined at which point in the runtime of the algorithm, the immediate results are sufficiently near to the final results to be useful. This is done because of the possibility to set the `uSVM_terminate` flag during the execution to force untimely termination.

### 5.1. Performance

We start this section by comparing the runtime of  $\mu$ SVM and Joachim's SVM<sup>light</sup> package on a personal computer for some example training sets. We also state the runtimes of  $\mu$ SVM on the ATmega16 microcontroller (section 4.2). We used four classes of example sets: EX A with  $n = 5$  and  $\ell = 5$ , EX B with  $n = 20$  and  $\ell = 20$ , EX C with  $n = 50$  and  $\ell = 10$ . For each of these classes, three randomly chosen example sets were tested. For all these tests, we took a linear kernel with error penalty  $C = 3.0$  and precision  $\varepsilon = 0.005$ . Additionally, the EX D.1 example set was chosen with the parameters  $n = 10$  and  $\ell = 30$ . The PC tests were performed on a PowerMac G5 with two 2 GHz processors and 2.5 GB RAM. The tests on the MCU were performed one time using an error cache (with EC) and one time without (w/o EC). The package SVM<sup>light</sup> was not evaluated on the MCU because of the (PC-oriented) big memory consumptions which prohibit execution on the target hardware.

The results in table 5.1 on page 31 can lead to the following conclusions:

- $\mu$ SVM performs quite well for small  $\ell$ , even for big  $n$ , but greatly loses performance with the growth of  $\ell$ .
- Growth of  $n$  affects operation with error cache much less than without error cache. This is because kernel evaluations are more expensive with big  $n$ .
- SVM<sup>light</sup> is more time efficient than  $\mu$ SVM on PC.

Next, we take a look at the numerical accuracy of  $\mu$ SVM. Therefore the results of  $\mu$ SVM on PC and SVM<sup>light</sup> are compared on the training example

Example set	SVM <sup>light</sup>	$\mu$ SVM on PC	$\mu$ SVM with EC	$\mu$ SVM w/o EC
EX A.1	0.010s	0.011s	2.7s	3.6s
EX A.2	0.009s	0.009s	2.4s	2.4s
EX A.3	0.010s	0.009s	3.6s	3.2s
EX B.1	0.017s	0.024s	88.2s	234.4s
EX B.2	0.019s	0.029s	170.0s	709.3s
EX B.3	0.018s	0.027s	196.0s	693.2s
EX C.1	0.011s	0.011s	10.8s	93.0s
EX C.2	0.012s	0.011s	13.1s	71.9s
EX C.3	0.012s	0.010s	9.9s	78.5s
EX D.1	0.121s	1.611s	> 40.0min	—

Table 5.1.: Runtime of SVM implementations

sets and summarized in table 5.2. Stated is the *norm* of the error vector  $v = (\lambda - \lambda^*, b - b^*)$  where  $(\lambda, b)$  is  $\mu$ SVM's solution and  $(\lambda^*, b^*)$  that of SVM<sup>light</sup>. Also specified is the relative error  $r = \|v\| / \|(\lambda^*, b^*)\|$ . SVM<sup>light</sup> was chosen as the numerical reference implementation because of its excellent reputation in the community.

Example set	absolute error	relative error
EX A.1	0.00023377	0.0481%
EX A.2	0.00004444	0.0040%
EX A.3	0.00008987	0.0102%
EX B.1	0.00136866	0.2153%
EX B.2	0.00069989	0.2498%
EX B.3	0.00147428	0.4782%
EX C.1	0.00058508	0.3161%
EX C.2	0.00668395	1.8747%
EX C.3	0.00255609	1.6248%
EX D.1	0.00285187	0.0273%

Table 5.2.: Numerical errors of  $\mu$ SVM

## 5.2. Speed of Convergence

We investigate how fast the solutions converge to the final result. We therefore measured the relative error of the intermediate results of the algorithm with respect to the final values. Some results are depicted in figures 5.1 and 5.2. The

other examples draw a similar picture. We see that after 20% of the execution time, the error is less than 10%.

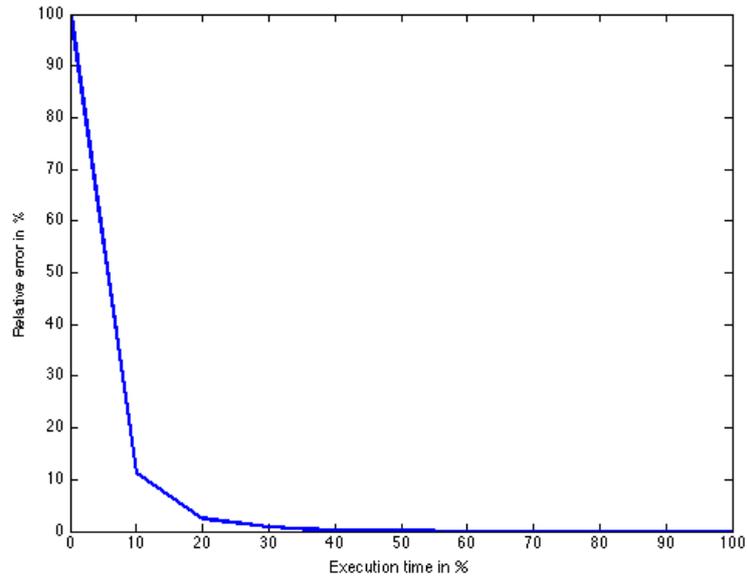


Figure 5.1.: Error progression of EX B.2

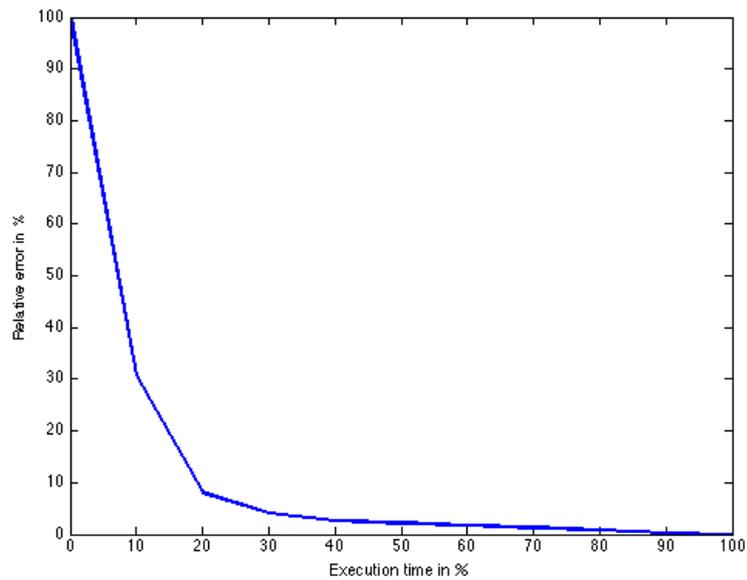


Figure 5.2.: Error progression of EX D.1

## 6. Summary

We started with a thorough introduction to convex analysis and derived the training process for Support Vector Machines. The Karush-Kuhn-Tucker conditions for solving constrained optimization problems, which are of great importance in practice, i. e., in the implementation, were explicitly written out in their general form and it was explained how these conditions can be applied to the SVM case.

An important part was generalizing the scalar product in  $\mathbb{R}^n$  to kernel functions, i. e., functions that are scalar products in some other Hilbert space. Mercer's theorem, a sufficient condition for a function to be a kernel function, was stated and proved in a very general setting ( $\sigma$ -finite measure spaces), which is probably a novelty in SVM literature.

We then introduced the SVM implementations of Osuna et al., Sequential Minimal Optimization (SMO) and SVM<sup>light</sup>. The SMO method was derived and investigated in more detail.

After illustrating the concepts of SVM training, we applied them in the implementation of  $\mu$ SVM. We have shown that despite the limited processing power and stringent memory space limitations in small microcontroller units, it is possible to use a fully-fledged SVM there. One problem is the long execution time of the implementation in the case of a large number of training examples. The experiments indicate, however, that it is often possible to stop the training process prematurely and still retain good numerical accuracy.

Future projects could focus on using  $\mu$ SVM in real-life applications or optimizing  $\mu$ SVM on other microcontroller units. Especially ones with more available memory space.

# Bibliography

- [Bur98] C. Burges. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2:121–167, 1998.
- [CHL00] C.-C. Chang, C.-W. Hsu, and C.-J. Lin. The analysis of decomposition methods for support vector machines. *IEEE Transactions on Neural Networks*, 11(4):1003, July 2000.
- [HUL93] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms I*. Springer, 1993.
- [Joa98] T. Joachims. Making large-scale support vector machine learning practical. In A. Smola B. Schölkopf, C. Burges, editor, *Advances in Kernel Methods: Support Vector Learning*. MIT Press, Cambridge, MA, 1998.
- [Mer09] J. Mercer. Functions of positive and negative type, and their connection with the theory of integral equations. *Philos. Trans. Roy. Soc. London*, 209:415–446, 1909.
- [OFG97] E. Osuna, R. Freund, and F. Girosi. Improved training algorithm for support vector machines. *Proceedings of the 1997 IEEE Workshop on Neural Networks for Signal Processing [1997] VII.*, pages 276–285, 24-26 Sep 1997.
- [Pla98] J. Platt. Fast training of support vector machines using sequential minimal optimization. In A. Smola B. Schölkopf, C. Burges, editor, *Advances in Kernel Methods: Support Vector Learning*. MIT Press, Cambridge, MA, 1998.
- [Vap98] V. N. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [Wei00] J. Weidmann. *Lineare Operatoren in Hilberträumen 1*. Teubner, 2000.

# A. Notation

$\mathbb{N}$	... set of positive integers
$\mathbb{R}$	... set of real numbers
$\overline{\mathbb{R}}$	... $\mathbb{R} \cup \{+\infty\}$
$B_r(x)$	... open ball with radius $r$ and center $x$ , i. e., $\{x \mid \ x\  < r\}$
$S^{n-1}$	... unit sphere in $\mathbb{R}^n$
$\inf$	... infimum
$\sup$	... supremum
$\lim a_i$	... limit of the net $(a_i)_i$
$\lim_{x \rightarrow x_0} f(x)$	... limit of the function $f$ at $x_0$
$\lim_{x \searrow x_0} f(x)$	... right-sided limit of the function $f$ at $x_0$
$\lim_{x \nearrow x_0} f(x)$	... left-sided limit of the function $f$ at $x_0$
$\nabla f(x)$	... gradient of $f$ at $x$
$f'(x)$	... derivative of $f$ at $x$
$f'_+(x)$	... right-sided derivative of $f$ at $x$
$f'_-(x)$	... left-sided derivative of $f$ at $x$
$D_v f(x)$	... directional derivative of $f$ at $x$ in direction $v$
$\partial f(x)$	... subdifferential of $f$ at $x$
$L^2(X)$	... space of quadratic integrable real-valued functions on $X$
$\ell^2(I)$	... space of quadratic summable real nets on $I$
$\langle x, y \rangle$	... inner product of $x$ and $y$
$ x $	... absolute value of $x$
$\ x\ $	... norm of $x$
$x^T$	... transpose of $x$
$e$	... $\sum 1/k!$
$\pi$	... $\Gamma(1/2)^2$
$\lambda$	... Lebesgue measure on $\mathbb{R}^n$
$\Gamma$	... gamma function, $x \mapsto \int_0^\infty t^{x-1} e^{-t} dt$
$\operatorname{sgn}$	... signum function
$\tanh$	... hyperbolic tangent function

## B. $\mu$ SVM Documentation

$\mu$ SVM is a Support Vector Machine (SVM) implementation for use on microcontroller units (MCUs). The two main functions are:

- `int uSVM_train()`
- `char uSVM_classify(float *x)`

Their usage is described in the following sections.

### B.1. Training

The training process is organized as follows.

1. Write the dimension  $n$  of the example vectors into `uSVM_n`.
2. Write the number  $\ell$  of example vectors into `uSVM_e11`.
3. Allocate memory for  $n \cdot (\ell + 1)$  vector components (either `char`, `int` or `float`, see B.3) and store the pointer in `uSVM_x`.
4. Allocate memory for  $\ell$  `char` variables and store the pointer in `uSVM_omega`.
5. Write example vectors with the `uSVM_WRITE(i,k,z)` macro. Here, the  $k$ th component of the  $i$ th vector is written to  $z$ .
6. Write labels of example vectors into the `uSVM_omega` array. The values here can only be `+1` and `-1`.
7. Select the kernel function by setting the `uSVM_ker` function pointer to one of the functions described in B.1.1.
8. Select the precision `uSVM_EPS` and the penalty paramter `uSVM.C`.
9. Call `uSVM_train()`.

The return value of `uSVM_train()` is either `-1` in case of an error and the number of support vectors otherwise. This number is also stored in the `uSVM_nSV` variable.

### B.1.1. Kernels

Currently available kernels are:

- `uSVM_scalar`: Linear kernel.
- `uSVM_gauss`: Gaussian kernel. The parameter  $\sigma$  can be modified with the `uSVM_GAUSS_SIGMA` macro. Default value is  $\sigma = 1$ .
- `uSVM_poly`: Polynomial kernel. The exponent  $p$  can be modified with the `uSVM_POLY_P` macro. Default value is  $p = 3$ .

### B.1.2. Early Termination

It is possible to terminate the training process ahead of time by setting the flag `uSVM_terminate`. This feature was added because experiments show that the computed values change only marginally after about 20% of the execution time.

## B.2. Classification

After training, new vectors can be classified by calling `uSVM_classify(x)`, where `x` is an array of  $n$  `float` values. The return value is either `+1` or `-1` reflecting the label that is given to the vector by the SVM.

## B.3. Compiler Flags

- `uSVM_X_INT` and `uSVM_X_FLOAT`: These flags select the data type used for training example vectors. Default is `char`.
- `uSVM_NO_ERROR_CACHE`: Disables the use of the error cache. The training algorithm needs less memory space, but is slower with this flag.
- `uSVM_GAUSS_SIGMA`, `uSVM_POLY_P`.

## B.4. Example

```
uSVM_ker = uSVM_scalar;  
uSVM_n = 5;
```

```
uSVM_ell = 5;
uSVM_C = 3.0;
uSVM_EPS = 0.005;
uSVM_x = malloc((uSVM_ell+1)*uSVM_n * sizeof(char));
uSVM_omega = malloc(uSVM_ell * sizeof(char));

uSVM_omega[0] = +1;
uSVM_x[0] = 2;
uSVM_x[1] = -3;
uSVM_x[2] = -1;
uSVM_x[3] = -5;
uSVM_x[4] = 2;
uSVM_omega[1] = -1;
uSVM_x[5] = 2;
uSVM_x[6] = -3;
uSVM_x[7] = 4;
uSVM_x[8] = -2;
uSVM_x[9] = -3;
uSVM_omega[2] = +1;
uSVM_x[10] = -1;
uSVM_x[11] = 0;
uSVM_x[12] = -1;
uSVM_x[13] = 4;
uSVM_x[14] = -5;
uSVM_omega[3] = +1;
uSVM_x[15] = -1;
uSVM_x[16] = 0;
uSVM_x[17] = -1;
uSVM_x[18] = 4;
uSVM_x[19] = -5;
uSVM_omega[4] = -1;
uSVM_x[20] = 1;
uSVM_x[21] = -1;
uSVM_x[22] = -1;
uSVM_x[23] = 4;
uSVM_x[24] = -1;

uSVM_train();

float *z = malloc(uSVM_n * sizeof(float));
z[0]=-1;
z[1]=0;
```

```
z[2]=-1;
```

```
z[3]=4;
```

```
z[4]=-5;
```

```
uSVM_classify(z);
```

```
free(z);
```